

RFP Attachment 4

FastForward R&D Draft Statement of Work

March 29, 2012



U.S. DEPARTMENT OF
ENERGY

Office of
Science

This work performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

Contents

1	INTRODUCTION	5
2	ORGANIZATIONAL OVERVIEW	6
2.1	The Department of Energy Office of Science.....	6
2.1.1	Advanced Scientific Computing Research Program	6
2.2	National Nuclear Security Administration	7
2.2.1	Advanced Simulation and Computing Program	7
3	MISSION DRIVERS	7
3.1	Office of Science Drivers.....	7
3.2	National Nuclear Security Administration Drivers	8
4	EXTREME-SCALE TECHNOLOGY CHALLENGES	8
4.1	Power Consumption and Energy Efficiency	8
4.2	Concurrency	9
4.3	Fault Tolerance and Resiliency	10
4.4	Memory and Storage Architecture	10
4.5	Programmability/Productivity.....	11
5	APPLICATION CHARACTERISTICS.....	12
6	ROLE OF CO-DESIGN	13
6.1	Overview	13
6.2	ASCR Co-Design Centers	14
6.3	ASC Co-Design Center.....	14
6.4	Proxy Apps.....	14
7	REQUIREMENTS	15
7.1	Description of Requirement Categories	15
7.2	Requirements for Research and Development Investment Areas	16
7.3	Common Mandatory Requirements	16
7.3.1	Solution Description (MR)	16
7.3.2	Research and Development Plan (MR)	16
7.3.3	Productization Strategy (MR)	17
7.3.4	Staffing/Partnering Plan (MR)	17
7.3.5	Project Management Methodology (MR)	17
7.3.6	Intellectual Property Plan (MR).....	17

8	EVALUATION CRITERIA	17
8.1	Evaluation Team	17
8.2	Evaluation Factors and Basis for Selection	17
8.3	Performance Features	18
8.4	Feasibility of Successful Performance	18
8.5	Supplier Attributes	19
8.5.1	Capability	19
8.6	Price of Proposed Research and Development	19
8.7	Alternate Proposals.....	19
	ATTACHMENT 1: PROCESSOR RESEARCH AND DEVELOPMENT REQUIREMENTS.....	20
	ATTACHMENT 2: MEMORY RESEARCH AND DEVELOPMENT REQUIREMENTS	25
	ATTACHMENT 3: STORAGE AND INPUT/OUTPUT RESEARCH AND DEVELOPMENT REQUIREMENTS	32

1 INTRODUCTION

The Department of Energy (DOE) has a long history of deploying leading-edge computing capability for science and national security. Going forward, DOE's compelling science, energy assurance, and national security needs will require a thousand-fold increase in usable computing power, delivered as quickly and energy-efficiently as possible. Those needs, and the ability of high performance computing (HPC) to address other critical problems of national interest, are described in reports from the ten DOE Scientific Grand Challenges Workshops¹ that were convened in the 2008–2010 timeframe. A common finding across these efforts is that scientific simulation and data analysis requirements are exceeding petascale capabilities and rapidly approaching the need for exascale computing. However, workshop participants also found that due to projected technology constraints, current approaches to HPC software and hardware design will not be sufficient to produce the required exascale capabilities.

In April 2011 a Memorandum of Understanding was signed between the DOE Office of Science (SC) and the DOE National Nuclear Security Administration Office (NNSA), Defense Programs, regarding the coordination of exascale computing activities across the two organizations. This led to the formation of a consortium that includes representation from seven DOE laboratories: Argonne National Laboratory, Lawrence Berkeley National Laboratory, Lawrence Livermore National Laboratory, Los Alamos National Laboratory, Oak Ridge National Laboratory, Pacific Northwest National Laboratory, and Sandia National Laboratories.

In July 2011, Argonne National Laboratory, on behalf of the seven aforementioned DOE labs, released a request for information (RFI) with the purpose of providing DOE SC and NNSA with information for planning the DOE Exascale Program. The RFI responses highlighted numerous challenges on the path to exascale and presented many innovative ideas to address those challenges. Funding for the DOE Exascale Program has not yet been secured, but DOE has compelling real-world challenges that will not be met by existing vendor roadmaps. Informed by responses to the exascale RFI, DOE SC and NNSA have identified three areas of strategic research and development (R&D) investment that will provide benefit to future extreme-scale applications:

- Processor technology
- Memory technology
- Storage and input/output (I/O)

These R&D activities will initially be pursued through a program called FastForward. The objective of the FastForward program is to initiate partnerships with multiple companies to accelerate the R&D of critical technologies needed for extreme-scale computing. It is recognized that the broader computing market will drive innovation in a direction that may not meet DOE's mission needs. Many DOE applications place extreme requirements on computations, data movement, and reliability. FastForward seeks to fund innovative new and/or accelerated R&D of technologies targeted for productization in the 5–10 year timeframe. The period of performance for any subcontract resulting from this request for proposal (RFP) will be two years. The consortium expects to establish an ongoing program to continue innovation in these and

¹ <http://science.energy.gov/ascr/news-and-resources/workshops-and-conferences/grand-challenges/>

additional technology areas. Contracts awarded through this RFP process may be eligible for additional funding to add work scope to accelerate further the critical technology R&D if Congress approves funding for this purpose.

The consortium is soliciting innovative R&D proposals in the areas of processor, memory, and storage, and I/O that will maximize energy and concurrency efficiency while increasing the performance, productivity, and reliability of key DOE extreme-scale applications. Due to the focus on extreme-scale applications, overall time to solution is also an important consideration. The goal is to begin addressing long-lead time items that will impact extreme-scale DOE systems later this decade. Technology roadmaps, as they exist today, threaten to have a hugely disruptive and costly impact on development of DOE applications and ultimately a negative impact on the productivity of DOE scientists.

Proposals submitted in response to this solicitation must address the impact of the proposed R&D on both DOE extreme-scale mission applications as well as the broader HPC community. Offerors are expected to leverage the DOE Co-Design Centers to ensure solutions are aligned with DOE needs. While DOE's extreme-scale computer requirements are a driving factor, these projects must also exhibit the potential for technology adoption by broader segments of the market outside of DOE supercomputer installations. This public-private partnership between industry and the DOE, initiated with FastForward, will aid the development of technology that reduces economic and manufacturing barriers to constructing exaflop-sustained systems, but also further DOE's goal that the selected technologies have the potential to impact low-power embedded, cloud/datacenter, and midrange HPC applications. This ensures that DOE's investment furthers a sustainable software/hardware ecosystem supported by applications across not only HPC but the broader IT industry. This will result in an increase in the consortium's ability to leverage commercial developments. It is not the consortium's intent to fund the engineering of near-term capabilities that are already on existing product roadmaps.

2 ORGANIZATIONAL OVERVIEW

2.1 The Department of Energy Office of Science

The SC is the lead Federal agency supporting fundamental scientific research for energy and the Nation's largest supporter of basic research in the physical sciences. The SC portfolio has two principal thrusts: direct support of scientific research and direct support of the development, construction, and operation of unique, open-access scientific user facilities. These activities have wide-reaching impact. SC supports research in all 50 States and the District of Columbia, at DOE laboratories, and at more than 300 universities and institutions of higher learning nationwide. The SC user facilities provide the Nation's researchers with state-of-the-art capabilities that are unmatched anywhere in the world.

2.1.1 Advanced Scientific Computing Research Program

Within SC, the mission of the Advanced Scientific Computing Research (ASCR) program is to discover, develop, and deploy computational and networking capabilities to analyze, model, simulate, and predict complex phenomena important to the DOE. A particular challenge of this program is fulfilling the science potential of emerging computing systems and other novel

computing architectures, which will require numerous significant modifications to today's tools and techniques to deliver on the promise of exascale science.

2.2 National Nuclear Security Administration

The NNSA is responsible for the management and security of the nation's nuclear weapons, nuclear non-proliferation, and naval reactor programs. It also responds to nuclear and radiological emergencies in the United States and abroad.

2.2.1 Advanced Simulation and Computing Program

Established in 1995, the Advanced Simulation and Computing (ASC) Program supports NNSA Stockpile Stewardship Programs' shift in emphasis from test-based confidence to simulation-based confidence. Under ASC, simulation and computing capabilities are developed to analyze and predict the performance, safety, and reliability of nuclear weapons and to certify their functionality. Modern simulations on powerful computing systems are key to supporting the U.S. national security mission. As the nuclear stockpile moves further from the nuclear test base through either the natural aging of today's stockpile or introduction of component modifications, the realism and accuracy of ASC simulations must further increase through development of improved physics models and methods requiring ever greater computational resources.

3 MISSION DRIVERS

3.1 Office of Science Drivers

DOE's strategic plan calls for promoting America's energy security through reliable, clean, and affordable energy, ensuring America's nuclear security, strengthening U.S. scientific discovery, economic competitiveness, and improving quality of life through innovations in science and technology. In support of these themes is DOE's goal to significantly advance simulation-based scientific discovery, which includes the objective to "provide computing resources at the petascale and beyond, network infrastructure, and tools to enable computational science and scientific collaboration." All the other research programs within the SC depend on the ASCR to provide the advanced facilities needed as the tools for computational scientists to conduct their studies.

Between 2008 and 2010, program offices within the DOE held a series of ten workshops² to identify critical scientific and national security grand challenges and to explore the impact exascale modeling and simulation computing will have on these challenges. The extreme scale workshops documented the need for integrated mission and science applications, systems software and tools, and computing platforms that can solve billions, if not trillions, of equations simultaneously. The platforms and applications must access and process huge amounts of data efficiently and run ensembles of simulations to help assess uncertainties in the results. New simulations capabilities, such as cloud-resolving earth system models and multi-scale materials models, can be effectively developed for and deployed on exascale systems. The petascale machines of today can perform some of these tasks in isolation or in scaled-down combinations

² <http://science.energy.gov/ascr/news-and-resources/workshops-and-conferences/grand-challenges/>

(for example, ensembles of smaller simulations). However, the computing goals of many scientific and engineering domains of national importance cannot be achieved without exascale (or greater) computing capability.

3.2 National Nuclear Security Administration Drivers

Maintaining the reliability, safety, and security of the nation’s nuclear deterrent without nuclear testing relies upon the use of complex computational simulations to assess the stockpile, to investigate basic weapons physics questions that cannot be investigated experimentally, and to provide the kind of information that was once gained from underground experiments. As weapon systems age and are refurbished, the state of systems in the enduring stockpile drifts from the state of weapons that were historically tested. In short, simulation is now used in lieu of testing as the integrating element. The historical reliance upon simulations of specific weapons systems tuned by calibration to historical tests will not be adequate to support the range of options and challenges anticipated by the mid-2020s, by which time the stewardship of the stockpile will need to rely on a science-based predictive capability.

To maintain the deterrent, the Nuclear Posture Review (NPR) insists that “the full range of Life Extension Program (LEP) approaches will be considered: refurbishment of existing warheads, reuse of nuclear components from different warheads, and replacement of nuclear components.” In addition, it is recognized that as the number of weapons in the stockpile is reduced, the reliability of the remaining weapons becomes more important. By the mid-2020s, the stewardship of the stockpile will need to rely on a science-based predictive capability to support the range of options with sufficient certainty as called for in the NPR. In particular, existing computational facilities and applications will be inadequate to meet the demands for the required technology maturation for weapons surety and life extension by the middle of the next decade. Evaluation of anticipated surety options is raising questions for which there are shortcomings in our existing scientific basis. Correcting those shortcomings will require simulation of more detailed physics to model material behavior at a more atomistic scale and to represent the state of the system. This pushes the need for computational capability into the exascale level.

4 EXTREME-SCALE TECHNOLOGY CHALLENGES

The HPC community has done extensive analysis³ of the challenges of delivering exascale-class computing. These challenges also apply more generally to extreme-scale HPC, regardless of whether or not the end result is an exaflop computer. In this section, we provide an overview of the most significant of these challenges.

4.1 Power Consumption and Energy Efficiency

All of the technical reports on exascale systems identify the power consumption of the computers as the single largest challenge going forward. Today, power costs for the largest petaflop systems

³ http://science.energy.gov/~media/ascr/ascac/pdf/reports/Exascale_subcommittee_report.pdf;
http://science.energy.gov/~media/ascr/pdf/program-documents/docs/Arch_tech_grand_challenges_report.pdf;
http://science.energy.gov/~media/ascr/pdf/program-documents/docs/Crosscutting_grand_challenges.pdf;
<http://www.cse.nd.edu/Reports/2008/TR-2008-13.pdf>; <http://www.exascale.org/mediawiki/images/2/20/IESP-roadmap.pdf>

are in the range of \$5–10M annually. To achieve an exascale system using current technology, the annual power cost to operate the system would be above \$2.5B per year with a power load of over a gigawatt (more than many power plants currently produce). To keep the operating costs of such a system in some kind of feasible range, a target of 20 megawatts has been established.

The power consumed by data movement will dominate the power budget of future systems. The power consumed in moving data between memory and processor is of particular concern. Historically a bandwidth/flop ratio of around 1 byte/flop has been considered a reasonable balance. With current double data rate three (DDR-3) memory technology, the energy cost to load a double-precision operand to memory is about 5×10^{-9} joules. For a current computer operating at 2 petaflop/s, the power required to maintain a 1 byte/flop ratio is about 1.25 MW. Extrapolating the JEDEC roadmap to 2020 and accounting for the expected improvements of DDR-5 technology, the total power consumption of the memory system would jump to 260 MW, well above the posited parameters for an exascale system. Even reducing the byte/flop ratio to 0.2—considered by some experts to be the minimum acceptable value for large-scale modeling and simulation problems—power consumption of the memory subsystem still would exceed 50 MW.

Achieving the power target for exascale systems is a significant research challenge. Even with optimistic expectations of current R&D activities, ***there is at least a factor of five gap between what we must have and what current research can provide.*** To get the additional factor of five improvements in power efficiency over projections, a number of technical areas in hardware design need to be explored. These may include: energy efficient hardware building blocks (central processing unit (CPU), memory, interconnect), novel cooling, and packaging, Si-Photonic communication, and power-aware runtime software and algorithms.

4.2 Concurrency

The end of increasing single compute node performance by increasing Instruction Level Parallelism (ILP) and/or higher clock rates has left explicit parallelism as the only mechanism in silicon to increase performance of a system. Scaling up in absolute performance will require scaling up the number of functional units accordingly, projected to be in the billions for exascale systems.

Efficiently exploiting this level of concurrency, particularly in terms of applications programs, is a challenge for which there currently are no good solutions. Memory latency further compounds the concurrency issue. We are already at or beyond our ability to find enough activities to keep hardware busy in classical architectures while long-time events such as memory references occur. While the flattening of clock rates has one positive effect in that such latencies will not get dramatically worse by themselves, the explosive growth in concurrency means that there will be substantially more of these high latency events; and the routing, buffering, and management of all these events will introduce even more delay. When applications then require any sort of synchronization or other interaction between different threads, the effect of this latency will be to exponentially increase the facilities needed to manage independent activities, which in turn forces up the level of concurrent operations that must be derived from an application to hide them.

Further complicating this is the explosive growth in the ratio of energy to transport data versus the energy to compute with it. At the exascale level, this transport energy becomes a front-and-

center issue in terms of architecture. Reducing the transport energy will require creative packaging, interconnect, and architecture changes to bring the data needed by a computation energy-wise “closer to” the function units. This closeness translates directly into reducing the latency of such accesses in creative ways that are significantly better than today's multi-level cache hierarchies.

4.3 Fault Tolerance and Resiliency

Resilience is a measure of the ability of a computing system and its applications to continue working in the presence of system degradations and failures. The resiliency of a computing system depends strongly on the number of components that it contains and the reliability of the individual components. Exascale systems will be composed of huge numbers of components constructed from VLSI devices that will not be as reliable as those in use today. It is projected that the mean time to interrupt (MTTI) for some components of an exascale system will be in the minutes or seconds range. Increasing evidence points to a rise in silent errors (faults that never get detected or get detected long after they generated erroneous results), causing havoc, which will only get more problematic as the number of components rises.

Exascale systems will continually experience failures, necessitating significant advances in the methods and tools for dealing with them. Achieving acceptable levels of resiliency in exascale systems will require improvement in hardware and software reliability, better understanding of the root cause of errors, better reliability, availability, and serviceability (RAS) collection and analysis, fault resilient algorithms and applications to assist the application developer, and local recovery and migration. The goal of research in this area is to improve the application MTTI by greater than 100 times, so that applications can run for many hours. Additional goals are to improve by a factor of 10 times the hardware reliability and improve by a factor of 10 times the local recovery from errors.

4.4 Memory and Storage Architecture

From a scientist's perspective, the ratio of memory to processor is critical in determining the size and type of problem that can be solved. Trends show a decrease in both memory size and bandwidth relative to system size. The rate of memory density improvement has gone from a 4-time improvement every three years to a 2-time improvement every three years (a 30-percent annual rate of improvement). Consequently, the cost of memory technology is not improving as rapidly as the cost of floating-point capability. Memory capacity is limited not only by fabrication and purchase cost; memory is also a large part of operational cost because it consumes a great deal of power. Without improvements in this area, it is anticipated that systems in the 2020 timeframe will suffer a 10-time loss in memory size relative to compute power.

Reduced memory bandwidth and increased latency will compound the memory capacity challenge. Neither bandwidth nor latency has improved at rates comparable to Moore's Law for processing units. On current petaflop systems, memory access at all levels is the limiting factor in most applications, so the situation for exaflop systems will be critical. Research in advanced memory technologies as well as optical interconnects and routers can provide critical improvements in latency and bandwidth. But work will also be needed in improving the efficiency of when and how data needs to move. Research options include better data analysis to anticipate needed data before it is requested (thus hiding latency), determining when data can be

efficiently recomputed instead of stored (reducing demands for bandwidth), and improved data layouts (to maximize the use of data when it is moved between levels.)

Regardless of efforts to reduce the amount of data that is moved and stored, exascale systems will still generate tremendous volumes of data that must be preserved. Systems ten years from now could have a billion cores, tens of petabytes of memory, and require hundreds of terabytes per second of I/O bandwidth to hundreds of petabytes of storage. This level of concurrency is well beyond the design point for today's HPC file systems. New approaches to application checkpointing, as well as alternative storage system paradigms, must be explored.

Furthermore, the current HPC storage stack relies heavily on low-cost, high-volume disk drives. Until recently, the number of disks required for capacity has been larger than the number required for bandwidth. However, disk capacity is now increasing much faster than disk performance and purchasing disks for bandwidth is cost-prohibitive. Solid-state drives (SSDs), while cost-effective for bandwidth, are cost-prohibitive for capacity. Future storage systems will no longer be able to assume an all disk storage device solution, and therefore, we anticipate solutions that involve hybrid storage or other technologies/concepts.

4.5 Programmability/Productivity

Programmability is the crosscutting property that reflects the ease by which application programs may be constructed. Programmability affects developer productivity and ultimately leads to the productivity of an HPC system as a tool to enable scientific research and discovery.

Programmability itself involves three stages of application development: (1) program algorithm capture and representation, (2) program correctness debugging, and (3) program performance optimization. All levels of the system, including the programming environment, the system software, and the system hardware architecture, affect programmability. The challenges to achieving programmability are myriad, related both to the representation of the user application algorithm and to underlying resource usage.

- **Parallelism**—sufficient parallelism must be exposed to maintain exascale operation and hide latencies. It is anticipated that 10-billion-way operation concurrency will be required.
- **Distributed Resource Allocation and Locality Management**—to make such systems programmable, the tension must be balanced between spreading the work among enough execution resources for parallel execution and co-locating tasks and data to minimize latency.
- **Latency Hiding**—intrinsic methods for overlapping communication with computation must be incorporated to avoid blocking of tasks and low utilization of computing resources.
- **Hardware Idiosyncrasies**—properties peculiar to specific computing resources such as memory hierarchies, instruction sets, and accelerators must be managed in a way that circumvents their negative impact while exploiting their potential opportunities without demanding explicit user control.

- Portability—application programs must be portable across machine types, machine scales, and machine generations. Performance sensitivity to small code perturbations should be minimized.
- Synchronization Bottlenecks—barriers and other over-constraining control methods must be replaced by lightweight synchronization overlapping phases of computation.

Novel architectures and execution models may increase programmability, thereby enhancing the productivity of DOE scientists.

5 APPLICATION CHARACTERISTICS

Multi-physics simulation is encountered in all missions supported by the DOE. "Multi-physics" numerical simulation is not simply simulation of complex phenomena on complex geometries. In its most simple form, multi-physics modeling involves two or more physical processes or phenomena that are coupled and that often require disparate methods of solution. For example, turbulent fluid simulations must be coupled to structural dynamics simulations, shock hydrodynamics simulations must be coupled to solid dynamics or radiation transport simulations, and atomic-level defects in electronic devices must be coupled to large-scale circuit simulations.

Computational modeling with multiple physics packages working together faces many challenging issues at the extreme scale. Among these are problems in which coupled physical processes have inherently different spatial and/or temporal attributes, leading to possibly conflicting discretizations of space and/or time, as well as problems where the solution spaces for the coupled physical processes are inherently distinct with some packages working in a real space while other parts of the solution require a higher dimensional solution space. As an example, for coupled radiation-hydrodynamics, the physical processes in the simulation impose inherently distinct demands on the computer architecture. Hydrodynamics is characterized by moderate floating-point computations with regular, structured communication. Monte Carlo particle transport is characterized by intense fixed-point computations with random communication. As a result, multi-physics simulations typically require well-balanced computer architectures in terms of processor speed, memory size, memory bandwidth, and interconnect bandwidth, at a minimum.

Typical simulations are composed of multiple physics packages, which advance a shared set of data throughout the problem simulation time. While the details vary among packages, all implementations require that multiple physics packages run concurrently. The algorithms developed to model these physics processes have disparate characteristics when implemented on parallel computer architectures. The data for the simulation is distributed across a mesh representing the phenomena modeled. For each element of this mesh, the algorithmic demands have been characterized in terms of memory requirements, communication patterns, and computational intensity described in the table below. These packages often have competing computation and communication requirements. Generally, the strategy is to compromise among the various competing needs of these packages, but an overall driving principle for major applications is to attain the maximum degree of accuracy in the minimum amount of time.

One key challenge of the algorithms used in multi-physics applications is a balance of the memory access characteristics where both the patterns and the size requirements differ

considerably and may fluctuate dramatically during the course of a calculation. Such variations impact both the communication patterns and the scaling characteristics of the codes. This is summarized in the following table:

Package	Memory per Mesh Element (KB)	Communication and Memory Access Patterns
A	0.2	Predictable with a modest amount of spatial and temporal locality
B	50–80	Predictable, but difficult to optimize, low spatial but high temporal locality
C	0.5–100	Unpredictable memory access, low spatial and low temporal locality
D	0.5	Predictable, with medium to high spatial and temporal locality

In addition to exascale computer architectures, multi-physics codes must also support other capacity-class computer architectures. Portability and high-level abstractions in the programming model will be critical. The complexity of the physics interaction in multi-physics codes tends to demand that the implementation have a single, shared code based on all computer architectures (that is, rewriting for boutique vendor hardware can quickly become a maintenance challenge). To date, mechanisms for expressing data hierarchies and optimization accessible by a given hardware realization have been closer to machine-level programming than high-level abstractions. As architectural complexities increase, research into appropriate abstractions in the programming model is needed. Additionally, improvements in the computational environment (for example, compilers and tools) are needed. This will become increasingly critical on exascale computer architectures. Addressing the issues of restrictions due to power constraints (how that impacts data layout), and heterogeneous node architectures are additional challenges⁴.

6 ROLE OF CO-DESIGN

6.1 Overview

The R&D funded through this RFP is expected to be the product of a co-design process. Co-design refers to a system design process where scientific problem requirements influence architecture design and technology, and architectural characteristics inform the formulation and design of algorithms and software. To ensure that future architectures are well-suited for DOE target applications and that DOE scientific problems can take advantage of the emerging computer architectures, major R&D centers of computational science are formally engaged in the hardware, software, numerical methods, algorithms, and applications co-design process.

Co-design methodology requires the combined expertise of vendors, hardware architects, system software developers, domain scientists, computer scientists, and applied mathematicians working together to make informed decisions about the design of hardware, software, and underlying algorithms. The future is rich with trade-offs, and give and take will be needed from both the

⁴ <http://www.sandia.gov/ascppc/ReferenceMaterials/Multiphysics%20white%20paper%20final.pdf>

hardware and software developers. Understanding and influencing these trade-offs is a principal co-design requirement.

ASCR and ASC are establishing multiple co-design centers that will be used as a vehicle to collaborate with vendors on R&D. The existing and planned co-design centers are at varying stages of deployment at the time of this RFP, and DOE is in the process of developing a strategy for issues such as cross-center collaboration, intellectual property protection, and overall governance. It is expected that much of this will be in-place and documented at the time FastForward awards are made.

6.2 ASCR Co-Design Centers

In mid-2011, ASCR granted the first three co-design awards, and additional ASCR centers may be established in the future. Each of these co-design centers is a distributed collaboration between multiple national laboratories and university partners. Each center has focused on a specific application that is an important driver for exascale and is using development of that application as a way to explore issues of mathematics, algorithms, computer science, systems software, and of course, hardware in the co-design process. For a detailed description of the ASCR co-design centers see <http://science.energy.gov/ascr/research/scidac/co-design/>.

6.3 ASC Co-Design Center

The NNSA labs and ASC program are defining a coordinated co-design strategy that leverages the work of the ASCR co-design centers while focusing on the unique needs of the ASC program. ASC is a mission-driven program with applications currently in use that are of importance to run at exascale in support of stockpile stewardship, namely the Engineering and Physics Integrated Codes (EPICs). To meet the key needs of the EPICs, ASC has established the National Security Applications (NSApp) Co-Design Center. NSApp will focus on these established applications as the drivers, and participate in co-design largely through proxy applications. Additional information is available at <https://asc.llnl.gov/codesign/>.

6.4 Proxy Apps

DOE will use proxy applications as the means to interact with our vendor partner(s) during the co-design process. These applications will be used both by the vendors to understand the effects of hardware tradeoffs, and also by integrated code team members and DOE researchers wishing to explore and develop new technologies, runtime systems, languages, programming models, algorithms, tools, file systems, and visualization techniques. Whenever possible, proxy apps are openly available—with occasional need to protect the original source under export-control rules or proprietary access rules in some cases where vendor modifications are supplied back to the co-design center.

Proxy apps can be grouped into three categories in increasing sophistication and fidelity to the actual applications (or packages) used in integrated design codes:

- **Micro-benchmarks:** Key tasks or processes extracted from numerical algorithms and solvers that allow for highly detailed study at the instruction level.

- **Kernels:** A combination of one or more fundamental primitives (the micro-benchmarks) integrated into a single executable most likely executing on a single type of device (for example, graphics processing unit (GPU) or multi-core CPU).
- **Skeleton apps:** Reproduced data flow of a simplified physics application with little or no attempt to investigate numerical performance. They are primarily useful in investigating network performance characteristics at large scale.
- **Mini apps:** These apps contain some of the dominant numerical kernels (or subsets thereof) contained in an actual application and produce simplifications of physical phenomena.
- **Compact apps:** These apps are representative of the actual application and may contain multiple kernels. In some cases compact apps may be full-fledged physics packages and as such are often restricted in their distribution.

ASC and ASCR co-design centers are in the process of developing and publishing their proxy apps. Some that are available today are:

TORCH	https://ftg.lbl.gov/projects/torch/
Mantevo	http://software.sandia.gov/mantevo
NERSC SSP	http://www.nersc.gov/research-and-development/performance-and-monitoring-tools/sustained-system-performance-ssp-benchmark/
LULESH	https://computation.llnl.gov/casc/ShockHydro/

7 REQUIREMENTS

7.1 Description of Requirement Categories

Requirements are either mandatory (designated MR) or target (designated TR-1, TR-2, or TR-3), and are defined as follows:

- MR are performance features essential to DOE requirements. An Offeror must satisfactorily address **all** MR to have its proposal considered responsive.
- TR, identified throughout this Statement of Work, are features, components, performance characteristics, or other properties that are important to DOE but will not result in a nonresponsive determination if omitted from a proposal. TR add value to a proposal and are prioritized by dash number. TR-1 is most desirable, while TR-2 is more desirable than TR-3.

TR-1s and MR are of equal value. The aggregate of MRs and TR-1s form a baseline solution. TR-2s are goals that boost a baseline solution, taken together as an aggregate of MRs, TR-1s, and TR-2s, into the moderately useful solution. TR-3s are stretch goals that boost a moderately useful solution, taken together as an aggregate of MRs, TR-1s, TR-2s, and TR-3s, into the highly useful solution.

7.2 Requirements for Research and Development Investment Areas

Detailed requirements for each of the three targeted R&D areas of investment are provided as Attachments to this document. A single proposal may address multiple areas of investment, that is, an Offeror need not submit a unique proposal for each area of investment on which it chooses to propose. Each proposal shall address all of the common MRs listed below. All of the MRs in each area of investment shall be included in the proposal.

7.3 Common Mandatory Requirements

The following items are mandatory for all proposals. That is, they must be present in any proposal for that proposal to be considered responsive and eligible for further evaluation.

7.3.1 Solution Description (MR)

Offeror shall describe the proposed R&D, with emphasis on how it will increase the performance of key DOE extreme-scale applications relative to energy usage while maintaining or increasing reliability and maintaining or decreasing runtimes.

Offerors shall discuss the innovative nature of the proposed R&D. Work that funds a company's current roadmap is not desired. Technology acceleration is acceptable if there is a clear DOE benefit and it is part of a broader strategy. The primary intent is to fund long-lead-time R&D objectives where significant advances can be made during the term of this program.

7.3.2 Research and Development Plan (MR)

Offeror shall provide a plan for conducting the proposed R&D, including timelines, milestones, and proposed deliverables. Deliverables shall be meaningful and measurable. Pricing shall be assigned to each milestone and deliverable. A schedule for periodic technical review by the DOE laboratories shall also be provided.

The R&D funded through this RFP is expected to be the product of a co-design process. More specifically, Offerors are expected to engage in co-design activities with DOE's ASC and ASCR Exascale Co-design Centers. The R&D plan shall include a discussion of how Offeror plans to collaborate with DOE researchers on co-design, with a detailed description of planned co-design efforts if known.

Some projects may develop a hardware prototype that demonstrates the value of the proposed concept. Others may perform a simulation or analysis that assesses the impact (or feasibility) of a proposed development. If funding provided through this RFP is insufficient to effectively demonstrate a concept or produce a prototype, Offerors shall provide a separate, non-binding budgetary estimate for follow-on work that would be needed to achieve this result. Do NOT include the estimated amount for this activity in the price for the R&D being proposed in response to this RFP. This follow-on work could be proposed in response to a future RFP, if one is issued.

We recognize that innovation involves risk. Proposals shall discuss technical and programmatic risk factors and the strategy to manage and to mitigate risk. If the planned R&D is not achieving the expected results, what alternatives will be considered? The amount of risk must be commensurate with the potential impact. Higher risk projects may be acceptable if the impact of the project is also high.

7.3.3 Productization Strategy (MR)

Offeror shall describe how the proposed technology will be commercialized, productized, or otherwise made available to customers. Offerors shall include identification of target customer base/market(s) for the technology. Offerors shall describe impact specifically on the HPC market as well as the potential for broad adoption. Solutions that have the potential for broader adoption beyond HPC are highly desired. Offerors shall indicate projected timeline for productization.

7.3.4 Staffing/Partnering Plan (MR)

Offerors shall describe staffing categories and levels for the proposed R&D activities. All lead and key personnel shall be identified by name and brief CVs for these personnel shall be provided. Any collaboration with other industry partners and/or universities shall be identified, and any key personnel from these partners/subcontractors shall be provided together with a description of their contributions to the overall effort.

7.3.5 Project Management Methodology (MR)

Project management and regular project status reporting are required. Offeror shall describe project management methodology and provide communication plan indicating method of communication (for example, written report, teleconference, and/or face-to-face meeting) and frequency (for example, weekly, monthly, and/or quarterly).

7.3.6 Intellectual Property Plan (MR)

Proposals shall include a plan for how each intellectual property (IP) item from each portion of the proposed R&D work will be handled, including requested IP ownership and licensing. Please consult RFP letter for information on Federal regulations concerning IP.

8 EVALUATION CRITERIA

8.1 Evaluation Team

The Evaluation Team includes representation from seven DOE laboratories: Argonne National Laboratory, Lawrence Berkeley National Laboratory, Lawrence Livermore National Laboratory, Los Alamos National Laboratory, Oak Ridge National Laboratory, Pacific Northwest National Laboratory, and Sandia National Laboratories, as well as Federal government representatives. Lawrence Livermore National Security (LLNS), as the entity awarding subcontracts as a result of this RFP, will act as the source selection official.

8.2 Evaluation Factors and Basis for Selection

Evaluation factors are mandatory requirements, performance features, supplier attributes, and price that the Evaluation Team will use to evaluate proposals. The Evaluation Team has identified the mandatory requirements, performance features and supplier attributes listed above and in each Attachment that should be discussed in the proposal. Offerors may identify and discuss other performance features and supplier attributes they believe may be of value to the Evaluation Team. If the Evaluation Team agrees, consideration may be given to them in the evaluation process. The Evaluation Team's assessment of each proposal's evaluation factors will

form the basis for selection. LLNS intends to select the responsive and responsible Offerors whose proposals contain the combination of price, performance features, and supplier attributes offering the best overall value to DOE. The Evaluation Team will determine the best overall value by comparing differences in performance features and supplier attributes offered with differences in price, striking the most advantageous balance between expected performance and the overall price. Offerors must, therefore, be persuasive in describing the value of their proposed performance features and supplier attributes in enhancing the likelihood of successful performance or otherwise best achieving the DOE's objectives for extreme scale computing.

LLNS desires to select two Offerors for each area of technology discussed in the Attachments to this SOW. However, LLNS reserves the right, based on the proposals received in response to the RFP, to select none, one, or more than two for any area of technology.

LLNS reserves its rights to: 1) make selections on the basis of initial proposals; 2) negotiate with any or all Offerors for any reason; and 3) award subcontract(s) based on a single proposal that addresses more than one Attachment area of technology.

8.3 Performance Features

The Evaluation Team will validate that an Offeror's proposal satisfies the MR. The Evaluation Team will assess how well an Offeror's proposal addresses the TR. An Offeror is not solely limited to discussion of these features. An Offeror may propose other features or attributes if the Offeror believes they may be of value. If the Evaluation Team agrees, consideration may be given to them in the evaluation process. In all cases, the Evaluation Team will assess the value of each proposal as submitted.

The Evaluation Team will evaluate the following performance features as proposed:

- How well the proposed solution meets the overall programmatic objectives expressed in the SOW
- The degree to which the technical proposal meets or exceeds any TR
- The degree of innovation in the proposed R&D activities
- The extent to which the proposed R&D achieves substantial gains over existing industry roadmaps and trends
- The extent to which the proposed R&D will impact HPC and the broader marketplace
- Credibility that the proposed R&D will achieve stated results
- Credibility of the productization plan for the proposed technology
- Realism and completeness of the project work breakdown structure

8.4 Feasibility of Successful Performance

The Evaluation Team will assess the likelihood that the Offeror's proposed research and development efforts can be meaningfully conducted and completed within the anticipated two-year subcontract period of performance. The Evaluation Team will also assess the risks, to both the Offeror and the DOE laboratories, associated with the proposed solution. The Evaluation

Team will evaluate how well the proposed approach aligns with the Offeror's corporate roadmap and the level of corporate commitment to the project.

8.5 Supplier Attributes

The Evaluation Team will assess the following supplier attributes.

8.5.1 Capability

The Evaluation Team will assess the following capability-related factors:

- The Offeror's experience and past performance engaging in similar R&D activities
- The Offeror's demonstrated ability to meet schedule and delivery promises
- The alignment of the proposal with the Offeror's product strategy
- The expertise and skill level of key Offeror personnel
- The contribution of the management plan and key personnel to successful and timely completion of the work

8.6 Price of Proposed Research and Development

The Evaluation Team will assess the following price-related factors:

- Reasonableness of the total proposed price in a competitive environment
- Proposed price compared to the perceived value
- Price tradeoffs and options embodied in the Offeror's proposal
- Financial considerations, such as price versus value

8.7 Alternate Proposals

An Offeror may submit an alternate proposal in the area of extreme scale computing technology. The Evaluation Team may evaluate alternate proposals for award consistent with the preceding information or as otherwise deemed necessary.

ATTACHMENT 1:

PROCESSOR RESEARCH AND DEVELOPMENT REQUIREMENTS

In this RFP, the term processor typically refers to the set of capabilities within a single microprocessor chip, or a tightly integrated set of processor capabilities that spans several chips (for example, chip stacks, chip carriers, chip sets, and other such approaches). Both architecture and process technologies are of interest. Key challenges include energy usage, performance, data movement, concurrency, reliability, and programmability, all of which are interrelated.

This processor R&D program focuses on innovations required to most effectively architect processors that scale well while reducing energy usage with high reliability. It is also acceptable to take advantage of decreasing feature sizes to replicate more copies of capabilities on a single integrated circuit to increase computational capability if it is in the context of other innovations.

A1-1 Key Challenges for Processor Technology

A1-1.1 Energy Utilization

Energy and power are key design constraints for processors. Techniques to minimize or constrain power used by computations while maintaining predictable behavior are needed. Possible areas include architectural features to improve application efficiency, advanced power gating techniques, near threshold operation, as well as packaging techniques such as 3D integration and silicon photonics to enable optics direct to a processor socket.

A1-1.2 Resilience and Reliability

Processor reliability is a critical concern, especially since future DOE supercomputers will utilize hundreds of thousands of processors. If FIT rates cannot be improved, the MTBI will fall to unacceptable levels. The ability to identify, contain, and overcome faults quickly is of paramount importance.

A1-1.3 On-Chip and Off-Chip Data Movement

Improved methods are needed for on-chip and off-chip data movement. The ability to move data efficiently limits the performance of many HPC applications. The energy required to move one bit of data within the processor and to memory must be reduced to a few picoJoules. In addition, improved memory interfaces can increase the effective bandwidth delivered to applications.

A1-1.4 Concurrency

Future increases in clock speeds are expected to be limited. As a consequence, processor companies are dramatically increasing concurrency (for example, more cores, greater instruction bundling, and multithreading) as feature sizes decrease. Managing this concurrency and the associated data movement is a considerable challenge. Many technologies could address the associated challenges in exploiting the available concurrency, including improved synchronization mechanisms, flexible atomic operations, and transactional memory.

Architectural mechanisms to handle work queue models efficiently could also improve application performance.

A1-1.5 Programmability

Achieving high performance on next-generation processors will be a challenge. Application developers will need to deal with massive concurrency and may need to manage locality, power, and resilience. A software ecosystem is needed to support the development of new applications and the migration of existing codes. Novel architectures and execution models may increase programmability and enhance the productivity of DOE scientists. Issues include the programmability of proposed architectures both in terms of complexity and the effort that will be required on the part of DOE scientists to achieve high performance.

A1-2 Areas of Interest

The following are examples of objectives and technologies that could be considered in processor R&D proposals that address DOE's extreme-scale computing needs. Some of the items below may only apply to certain architectures, and some may be mutually exclusive. ***Proposals are not limited to these areas, and alternative topics are encouraged.***

A1-2.1 Energy Utilization

- Advances that improve the power efficiency of processors
- Advances in measurement and application control of power utilization
- Advances that support high-performance, power-efficient processor integration with memory, optics, and networking
- Techniques to reduce cooling energy requirements

A1-2.2 Resilience and Reliability

- Advances that improve the resiliency or reliability of processors, for example, improved fault detection and correction on chip
- Advances that permit automatic rollback (within a window) after a fault or synchronization error
- Advances that demonstrate hardware/software resilience tradeoffs to improve overall time to solution

A1-2.3 On-Chip and Off-Chip Data Movement

- Advances that allow extremely low-latency response to incoming messages
- Improvements to the performance and energy efficiency of messaging, remote memory access, and collective operations

- Advances that allow explicit (software controlled) movement of data in and out of various on-chip memory regions (for example, levels of cache)
- Hardware support for active messages
- Hardware support for large numbers of short messages to achieve low latency
- Integration of a NIC on the processor
- Other hardware mechanisms for eliminating overhead

A1-2.4 Concurrency

- Advances that improve the scalability of processor designs as the number of processing units per chip increase
- Advances that address the inherent scaling and concurrency limits in applications
- Advances that improve the efficiency of process or thread creation and their management
- Advances that reduce the synchronization and activation time of large numbers of on-chip threads
- Advances in transactional and speculative processing techniques that significantly aid DOE applications

A1.2.5 Programmability

- Advances that significantly improve the performance and energy efficiency of arithmetic patterns common to DOE applications but are not well supported by today's processors, for example, short vector operations such as processing in vector registers
- Advances that allow efficient computation on irregular data structures (for example, compressed sparse matrices and graphs)
- Research to determine the most effective option(s) for cache and memory coherency policies; configurable coherency policy and configurable coherence or NUMA domains may be options; coherency policy might also be a power management tool
- Support for a global address space

A1-3 Performance Metrics

Offeror shall estimate or quantify the impact of the proposed technology over industry roadmaps and trends. This information shall be provided for all of the metrics listed below. If Offeror determines that a particular metric is not applicable to the technology being proposed, then Offeror shall explain why they believe the metric is not relevant and shall replace that metric with an alternate *meaningful* metric.

Quantities specified should reflect solutions that are productized in the 2020 timeframe. These metrics are independent, but a solution that can deliver advances in more than one metric is more desirable than one that solves only one metric at the expense of the others. The most meritorious

improvements will make substantial gains over industry roadmaps/trends and substantiate a convincing path to achieving the extreme-scale technology characteristics required by DOE.

- Node and socket power requirements
- Processor computational capability per watt
- FIT rate per socket
- Error detection, correction, and coverage of hard and soft error types
- Energy per bit for data transfers
- Computational capacity per node
- Data motion overhead (and avoided data motion)
- Aggregate bandwidth delivered to memory

A1-4 Target Requirements

The requirements below apply to supercomputers that will be deployed at the end of this decade to meet DOE mission needs. As previously stated, Offerors need not address all problem areas, and thus the Offeror need not respond to TR below if the proposed capability does not address that problem area. In all TR responses that are provided, Offeror should discuss what progress will be made in the next two years and describe what follow-on efforts will be needed to fully achieve these goals. Offeror should describe in detail how the metric will be evaluated, including the measurement method that will be used (for example, simulation or prototype) and any assumptions that will be made.

In the discussion below, a node is defined as the smallest physical unit of hardware that contains a processor chip(s), memory, and at least one network connection to connect to other such units.

A1-4.1 Energy Utilization (TR-1)

An energy and concurrency efficient processor that achieves high performance on a broad range of DOE applications (for example, the co-design center applications described previously) is highly desired. Solutions should realize greater than 50 GF/Watt at system scale while maintaining or improving system reliability.

A1-4.2 Resilience and Reliability (TR-1)

Mean Time to Application Failure (TR-1). Processor designs should make advances that lead to a mean time to application failure requiring user or administrator action of six days or greater in a 2020 exascale system, as determined by estimates of system component FIT rates and application recovery rates.

Wall-Time Overhead (TR-1). The wall-time overhead to handle automatic fault recovery should not reduce application performance by more than half.

A1-4.3 On-Chip Data Movement (TR-2)

The aggregate node memory bandwidth should be at least 4 TB/s to a greater than 100-GB region of memory. It is highly desirable for a node to have 320–640 GB of memory. If a hierarchical memory structure is needed, then the proposal should discuss how the hierarchy and the bandwidth between each tier will be managed.

A1-4.4 Off-Chip Data Movement (TR-2)

The total bandwidth between a node and the interconnect should be greater than 400 GB/s.

A1-4.5 Concurrency (TR-2)

To keep system sizes manageable, the overall performance of a node should be greater than 10 TF.

A1-4.6 Programmability (TR-1)

Solutions will need a software ecosystem that supports the development of new applications, the migration of existing applications, application maintenance, and application portability, while enabling DOE scientists to achieve high performance with no more effort than is required for today's high-end computers. Offeror should describe in detail how this will be accomplished.

ATTACHMENT 2:

MEMORY RESEARCH AND DEVELOPMENT REQUIREMENTS

Current roadmaps to develop future commodity memory components are not on track to meet DOE requirements for HPC systems, and there appears to be little industry consensus on a way forward that will realistically enable memory systems to meet size and performance goals for systems at server levels or higher while living within stringent power limits. The power consumed by data movement will dominate the power consumption profile of future systems. Chief among these concerns is the power consumed by memory technology. Failure to address these concerns with early technology investment will force DOE to accept undesirable trade-offs regarding power consumption and breadth of applications that can run effectively on the system.

A2-1 Key Challenges for Memory Technology

Following are some areas of emphasis in memory technology based on the requirements of DOE's application workload. None of these need to be construed as pointing to specific prescribed solutions.

A2-1.1 Energy Consumption

Power consumption is a leading design constraint for future systems. Chief among these concerns is the power consumed by memory technology, which would easily dominate the overall power consumption of future systems if we continue along current technology trends. The target for an exaflop system in 2020 is 20 megawatts for the complete system. If we extrapolate commodity DDR memory technology trends out to 2020, the memory system alone would eclipse the target power budget and make future HPC systems of all scales less effective. *FastForward would like to develop memory technologies to improve the energy efficiency of memory while improving capacity, bandwidth, and resilience.*

A2-1.2 Memory Bandwidth

Memory bandwidth has always been a major bottleneck for the performance of HPC applications. As core count of processors have increased, the memory bandwidth available to each core has significantly decreased. Higher memory bandwidth enables a wider array of algorithms to fully utilize available computing performance. We recognize that improvements in memory bandwidth must be balanced against power consumption and capital costs incurred. *FastForward will emphasize the development and acceleration of technology to increase memory bandwidth while keeping cost, reliability, and power consumption under control.*

A2-1.3 Memory Capacity

The rate of improvement for DRAM density has slowed in recent year from quadrupling every three years to doubling every three years. In comparison, logic density and the cost of flops is improving at a much more rapid rate. The consequence is that we expect lower memory capacity per peak computational performance than in past machines. This is of concern for DOE applications because increased problem resolution requiring larger memory capacity is at least as

important for many scientific applications as improvements in computational performance. Worse yet, technology roadmaps out to 2020 indicate we can get high-capacity memory with low bandwidth and low-capacity memory with high bandwidth—but not both. However, the DOE mission need and scientific objectives require improvements in both increased problem sizes (limited by memory capacity) and performance (limited by memory bandwidth) in the same memory space. A solution that delivers one or the other (but not both), will fail to meet mission objectives. *The FastForward program is interested in accelerating and developing new technology options that can deliver both capabilities (bandwidth and capacity) in the same cost-effective package.*

A2-1.4 Reliability

Components that are otherwise reliable in consumer applications that contain only a handful of devices have high aggregate failure rates for scalable HPC systems that typically include millions of components. Even in today's HPC systems and large-scale datacenters, memory DIMMS are among the most common sources of hardware failure. A recent large-scale field study by Google and the University of Toronto has shown that DRAM failure rates are much higher than originally anticipated⁵. *For scalable systems, FastForward is interested in developing and accelerating technologies that dramatically reduce DRAM component failure rates over a baseline that is largely set by smaller scale consumer devices.*

A2-1.5 Error Detection/Correction/Reporting

With respect to component failure rates (reliability), there are concerns about the ability of modern error detection and correction technology to keep up with the increased rate of transient errors. For scalable HPC systems and large-scale data centers, there is an increased observation of uncorrectable errors (double-bit or burst errors). Even more worrisome are the increased incidence of silent errors, which are already apparent in modern HPC systems. It would greatly improve the usability of these resilience features if more comprehensive error detection and reporting technology were available (for example, S.M.A.R.T. technology for system boards). *FastForward is seeking technologies to improve and even scale our ability to detect and correct transient errors, and to greatly reduce the possibility of silent errors in large-scale systems.*

A2-1.6 Processing in Memory

An alternative approach to improving effective memory bandwidth is to embed computing operations within the memory component to reduce pressure on memory bandwidth. At minimum, this includes embedding basic element/word-granularity operations such as atomic memory operations and synchronization primitives in the memory to eliminate round-trips of data movement between the processor and the memory. At a medium level of integration, one could embed vector-primitives such as strided gather operations, general gather/scatter, and checksum operations (for end-to-end error detection) in the memory system to reorganize areas of memory to improve data transfer performance. General-purpose processing-in-memory is the most extreme and general approach to embedding a processing capability into the memory

⁵ B. Schroeder, E. Pinheiro, W-D. Weber, "DRAM Errors in the Wild: A Large-Scale Field Study," SIGMETRICS/Performance'09, ACM, Seattle WA. 2009.

subsystem. *FastForward is interested in novel ideas for embedding processing in memory to improve data transfer efficiency or even eliminate the need to move data off of the memory chip.*

A2-1.7 Integration of NVRAM Technology

Solid-state storage technology (FLASH and other forms of NVRAM) has found a way into consumer and HPC systems primarily as disk/file system technology. However, we see many opportunities for improved performance and capability if NVRAM were integrated directly into the memory hierarchy rather than as a disk replacement. For example, for checkpointing/resilience technology, node-local NVRAM could offer substantial benefits if it can be made substantially more trustworthy and reliable than the active DRAM memory from which the state is being “checkpointed.” On-chip NVRAM could preserve local register or pipeline state to support microcheckpointing for resilience or instant power-down operation for chips (which are useful in the consumer space, too). NV memory can be used to hold tables and data items that are seldom written to relieve some pressure from the DRAM portion of the memory system. NVRAM-backed DRAM could enable power-off of areas of memory that are un-used or under-utilized. *FastForward is seeking novel applications and solutions involving deeper integration of NVRAM technology in the memory hierarchy.*

A2-1.8 Ease of Programmability

As novel technologies are added to computer systems, the memory hierarchies can become very complex and require management from the application. An addition such as high-speed scratchpad memory or software-managed caches creates disparate memory spaces with varying performance characteristics and capacities. It is desirable to support a broader ecosystem of software for the device that will ensure the features will continue to be supported across systems. *Fast Forward is seeking novel hardware and software solutions to simplify the management of deep memory hierarchies.*

A2-2 Areas of Interest

Below are some areas of technology development and acceleration that could be considered in memory R&D proposals to address DOE’s extreme-scale computing needs. ***Proposals are not limited to these areas, and alternative topics are encouraged.***

- Technologies to improve the energy efficiency of memory while improving capacity, bandwidth, and resilience
- Technology to increase memory bandwidth while keeping cost, reliability, and power consumption under control
- New technology options that can deliver both bandwidth and capacity in the same cost-effective package
- Technologies that dramatically reduce DRAM and NVRAM component failure rates over a baseline that is largely set by smaller scale consumer devices
- Technologies to improve and scale the ability to detect and correct transient errors, and to prevent the incidence of silent errors in large-scale systems

- Novel ideas for embedding processing in memory to improve data transfer efficiency or even eliminate the need to move data off the memory chip
- Novel applications and solutions involving deeper integration of NVRAM technology in the memory hierarchy
- Novel hardware and software solutions to simplify the management of deep memory hierarchies

A2-3 Performance Metrics (MR)

Offeror shall estimate or quantify the impact of the proposed technology over industry roadmaps and trends. This information shall be provided for all of the metrics listed below. If Offeror determines that a particular metric is not applicable to the technology being proposed, then Offeror shall explain why they believe the metric is not relevant and shall replace that metric with an alternate *meaningful* metric.

Quantities specified should reflect solutions that are productized in the 2020 timeframe. These metrics are independent, but a solution that can deliver advances in more than one metric is more desirable than one that solves only one metric at the expense of the others. The most meritorious improvements will make substantial gains over industry roadmaps/trends and substantiate a convincing path to achieving the extreme-scale technology characteristics required by DOE.

A2-3.1 DRAM Performance Metrics

Energy per Bit. This is defined as the energy needed to completely run memory, counted per bit of data moved, including a short length of interconnect (~2 cm) and the end-points (the complete memory chip, SerDes, wire losses, and memory controller on the CPU side). Offeror shall specify projected energy per bit for proposed DRAM solutions. Offeror shall describe any assumptions used in calculating this metric and how it will be measured. Seven picojoules per bit is considered the baseline value for this metric.

Aggregate Bandwidth per Socket (DRAM or Suitable Replacement for DRAM). This is defined as the data bandwidth delivered to a processor chip comprising the “socket.” A socket is defined as the smallest physical unit of hardware that contains one processor chip, memory, and at least one network connection to connect to other such units. Offeror shall specify both the peak performance as well as what *measured* performance can be expected for different access patterns, and how bandwidth would be measured. One TB/s is considered the baseline value for this metric.

Memory Capacity per Socket. This is defined as the *usable* data capacity per socket. Offeror shall specify the projected DRAM capacity and how it relates to other memory metrics such as bandwidth. One hundred GB is considered the baseline value for this metric.

FIT Rate per Node. This is the total soft-error FIT rate for the portion or fraction of a memory system, per node. A node is defined as the smallest physical unit of hardware that contains processor chip(s), memory, and at least one network connection to connect to other such units. The FIT rate is defined as the number of unrecoverable soft errors per billion hours of operation. This is not the sum of FIT rates but assumes additional error detection and recovery, for

example, possibly with spare components. Offeror shall describe how the FIT rate will be measured, the cost of recovery from transient errors (time/power), and assumptions used in the fault model. A FIT rate of less than 1000 is considered the baseline value for this metric.

Error Detection. Offeror shall describe technologies that will significantly improve error detection, recovery, and reporting. Offeror shall describe in detail tests that would demonstrate how error detection coverage, reporting, and recovery have been improved over the baseline. ECC + bit steering is considered the baseline for this metric.

Processing in Memory. Offeror shall describe the degree to which any proposed processing in memory technology will reduce data movement in target DOE codes. Offeror shall describe the programming model that will make these features productive for software developers. At a minimum, solutions must include support for atomics in memory.

Programmability/Usability. Offeror shall describe how any proposed memory technology feature would be integrated into a productive programming environment. Offeror shall specify projected improvements in productivity of end users and software developers. At a minimum, solutions must make existing programming models easier to use.

A2-3.2 NVRAM Performance Metrics

NVRAM Integration. Offeror shall describe the cell technology and architecture for NVRAM integration, and at what level of the node architecture this NVRAM would be integrated (for example, tightly integrated devices such as NVRAM-backed register files within a CPU versus loosely integrated SSD-like devices for node-level data storage).

Energy per Bit. This metric is largely the same as the DRAM energy per bit. However, the manner for calculating the energy will be highly dependent on where the NVRAM is integrated into the system. For NVRAM power, we are more interested in the energy required to write data versus read data in the proposed cell rather than the cost of data movement. Offeror shall specify projected energy per bit for proposed NVRAM solutions. Offeror shall describe all assumptions and specific tests that would be used to assess this energy metric. Offeror shall explain how the energy per bit and performance relates to wear-out rates for storage cells, if applicable to the proposed NVRAM technology.

Aggregate Bandwidth per Socket. This is defined as the data bandwidth delivered to a processor chips comprising the “socket.” A socket is defined as the smallest physical unit of hardware that contains one processor chip, memory, and at least one network connection to connect to other such units. Offeror shall specify both the peak performance for NVRAM as well as what *measured* performance can be expected for different access patterns, and how bandwidth would be measured.

Capacity per Socket. This is defined as the *usable* data capacity per socket. Offeror shall specify the projected NVRAM capacity. Five hundred GB is considered the baseline for this metric.

FIT Rate per Node. This is the total soft-error FIT rate for the portion or fraction of a memory system, per node. A node is defined as the smallest physical unit of hardware that contains processor chip(s), memory, and at least one network connection to connect to other such units. The FIT rate is defined as the number of unrecoverable soft errors per billion hours of operation. This is not the sum of FIT rates but assumes additional error detection and recovery, for

example, possibly with spare components. Offeror shall describe how the FIT rate would be measured, the cost of recovery from transient errors (time/power), and what the assumptions are that feed in their fault model. We are particularly interested in how NVRAM technologies can be made substantially less prone to failure so that they can be used as a reliable backing store to recover from errors/faults at the node level.

Error Detection. Offeror shall describe technologies that can significantly improve NVRAM error detection, recovery, and reporting. Offeror shall describe in detail tests that would demonstrate how error detection coverage, reporting, and recovery have been improved over the baseline.

Programmability/Usability. Offeror shall describe how any proposed NVRAM memory technology feature would be integrated into a productive programming environment. Offeror shall specify projected improvements in productivity of end users and software developers.

A2-4 Target Requirements

The requirements below apply to supercomputers that will be deployed at the end of this decade to meet DOE mission needs. As previously stated, Offerors need not address all problem areas, and thus the Offeror need not respond to TR below if the proposed capability does not address that problem area. In all TR responses that are provided, Offeror should discuss what progress will be made in the next two years and describe what follow-on efforts will be needed to fully achieve these goals. Offeror should describe in detail how the metric will be evaluated, including the measurement method that will be used (for example, simulation or prototype) and any assumptions that will be made.

A2-4.1 Energy per Bit

Reduced Energy per Bit (TR-1)

Energy per bit should be 5 picojoules or less end-to-end. End-to-end is defined as including full path from memory to register on processor chip, including the memory component and cost of accessing the memory cell in the memory component.

Greatly Reduced Energy per Bit (TR-2)

Energy per bit should be 1 picojoule or less end-to-end.

A2-4.2 Aggregate Delivered DRAM Bandwidth

Improved Aggregate Delivered DRAM Bandwidth Per Socket (TR-1)

Aggregate delivered bandwidth per socket for DRAM or equivalent should be 4 TB/s or greater over a distance of 5 cm or more.

Greatly Improved Aggregate Delivered DRAM Bandwidth Per Socket (TR-2)

Aggregate delivered bandwidth per socket for DRAM or equivalent should be 10 TB/s or greater over a distance of 5cm or more.

A2-4.3 Memory Capacity per Socket

Increased DRAM Capacity per Socket (TR-1)

Memory capacity per socket for DRAM or equivalent should be 500 GB or greater with preference for “fast” memory per item 4.2 above.

Greatly Increased DRAM Capacity per Socket (TR-2)

Memory capacity per socket for DRAM or equivalent should be 1 TB or greater with preference for “fast” memory per item 4.2 above.

Increased NVRAM Capacity per Socket (TR-1)

Memory capacity per socket for NVRAM or equivalent should be 1 TB or greater with preference for greatly improved reliability per items 4.4 and 4.5 below.

A2-4.4 FIT Rate per Node

Improved FIT Rate per Node (TR-1)

FIT rate per node should not exceed 100.

Greatly Improved FIT Rate per Node (TR-2)

FIT rate per node shall not exceed 10.

A2-4.5 Error Detection Coverage and Reporting

Reduction in Silent Errors (TR-1)

Solution should propose and estimate ways to greatly reduce possible rates of silent errors.

End-to-End Error Detection and Recovery (TR-2)

Solution should provide complete end-to-end error detection and recovery, including data paths.

A2-4.6 Advanced Processing in Memory Capabilities

Vector Operations and/or Gather/Scatter (TR-1)

Processing in memory solutions should include vector operations and/or gather/scatter.

General Purpose Processor in Memory (TR-2)

Offeror should implement a general-purpose processor-in-memory solution.

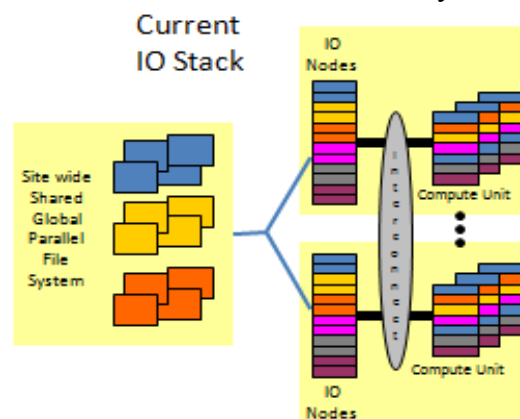
A2-4.7 Enhanced Programmability/Usability (TR-1)

Offeror should include full language support for new memory features.

ATTACHMENT 3: STORAGE AND INPUT/OUTPUT RESEARCH AND DEVELOPMENT REQUIREMENTS

A3-1 Key Challenges for Storage and Input/Output Technologies

In the petascale HPC era and before, the storage stack used by the extreme scale HPC community is fairly homogeneous across sites. On the compute edge of the stack, file system clients or I/O forwarding services direct I/O over an interconnect network to a relatively small set of I/O nodes. These nodes forward the requests over a secondary storage network to a disk-based parallel file system. Sizes of the systems at the petascale are hundreds of thousands of cores with hundreds of pebibytes (PiB) of memory. Disk systems at this scale involve thousands of disk devices and run at hundreds of gigabytes (GiB)/s with PiB of storage.



The current I/O stack described in the text.

The above description is provided to help explain the current storage and I/O deployments. This solicitation is not implying that proposed solutions need to fit into this model. For example, just because parallel file systems are in use now, Offerors should not assume this will be true in the future. Furthermore, given the innovative nature of this solicitation, we desire the focus be on new solutions beyond current product roadmaps.

Systems five years from now may be tens to hundreds of thousands of nodes in size, incorporating PiB of memory, and leveraging a low-latency network providing tens of GiB per second of bandwidth, per link. Systems ten years from now could have a billion cores, tens of PiB of memory, and require tens to hundreds of TiB/s of I/O bandwidth with hundreds of PiB of storage. A wide range of network configurations are possible, including but not limited to “fat” trees, dragonflies, tori, and hybrids. In all cases, such networks will almost certainly support hardware-assisted operations such as remote direct memory access and, perhaps, network operations in support of concurrency, such as semaphores or other atomics.

The following table outlines the expected extreme computing environment:

Year	2012	2015	2018	2020
Nodes	10–100 K	10–100 K	10 K–1 M	100 K–1 M
System Wide Concurrency	100 K–1 M	1–10 M	10 M–100 M	100 M–1 B
Memory (application byte addressable DRAM)	1–4 PiB	4–10 PiB	10–30 PiB	30–60 PiB
Scratch Size	10–100 PiB	50–300 PiB	200–900 PiB	600–3000 PiB

Year	2012	2015	2018	2020
Mean Time To Application Interrupt (without check-pointing or other resilience mechanism)	1–5 Days	1 Day	12 Hours	6 Hours
Time to Dump Memory	<2000 s	<1000 s	<600 s	<300 s
Peak I/O Burst Rates	2 TiB/s	10 TiB/s	50 TiB/s	200 TiB/s
Metadata Transaction Rates	100K/sec	1M/sec	10M/sec	100M/sec
Storage System Mean Time to Application Visible Interrupt	20 Days	18 Days	16 Days	14 Days

DOE desires solutions for a general-purpose supercomputer but will also consider technologies and I/O system or software architectures for special-purpose machines. This description is not meant to limit a proposal, only to communicate our best guess about the future.

I/O workloads consist of all of the following:

- Defensive output (bulk synchronous checkpointing—dumping of large portions of core memory from many cores simultaneously), from non-contiguous parts of main memory scattered over a billion cores
- Output of results, both large scale (billions of cores) and small scale (one core), from non-contiguous parts of main memory scattered over a billion cores
- Input of small configuration data to be sent to non-contiguous parts of main memory scattered over a billion cores
- Input of large amounts of data for analysis or restarting sent to non-contiguous parts of main memory scattered over a billion cores
- Input of dynamically loaded libraries, run-time linking of shared libraries, executables, and other system demands for I/O to service a million node, billion core-class environment

In addition to understanding the current and anticipated future environments and workloads, it is important to understand the economics associated with providing a storage and I/O service for extreme scale HPC environments. The current storage stack used in HPC is threatened by trends in disk technology and harsh economic realities. Until recently, the number of disks required for capacity in extreme HPC environments has been larger than the number required for bandwidth. In other words, buying the number of disks required for capacity has provided excess bandwidth essentially for free. However, disk capacity is increasing much faster than disk performance. New technologies such as shingled disks are only exacerbating this trend. The result is that the number of disks required for capacity has now become fewer than the number required for bandwidth. Unfortunately, purchasing disks for bandwidth is cost-prohibitive. Solid-state drives (SSDs), however, are cost-effective for bandwidth but cost-prohibitive for capacity.

Consequently we believe the extreme scale HPC storage environments of the future will no longer be an all disk device solution, and therefore other solutions will be required.

This solicitation seeks fundamental solutions and technologies with capabilities superior to today's offerings. Offerors may propose anything from an entire replacement of the I/O stack for extreme scale HPC to replacement or greatly enhanced versions of existing solutions. R&D well beyond current product roadmaps is desired. The Offeror may propose R&D to address any number of the challenges previously discussed as well as problems that are not stated that the proposer believes will exist in the future. The Offeror should propose R&D efforts and possibly prototypes to solve extreme-scale HPC storage and I/O problems. Proposed R&D should take into consideration the extreme-scale HPC environments of the future, anticipated workloads, and economic realities.

A3-2 Areas of Interest

Below are some areas of technology development and acceleration that could be considered in storage and I/O R&D proposals to address DOE's extreme-scale computing needs. ***Proposals are not limited to these areas, and alternative topics are encouraged.***

A3-2.1 Reliability/Availability/Manageability

Storage subsystems for today's extreme-scale HPC environments largely are made up of global parallel file systems that are disk based. These software systems were architected at least a decade ago and have fallen short of meeting our needs. It has become clear that current storage systems are not as reliable, available, nor as easily managed as befits a production service for our supercomputer centers. Looking forward, storage and I/O systems should assume failed components and infrastructure as the norm. Many existing storage and I/O systems simply cannot operate effectively in such an environment. Resilient operation in the presence of compromised infrastructure could be fundamental to the software design of the entire I/O stack. Availability/MTTI targets are provided in the environment table above.

A3-2.2 Metadata

It is expected that any storage and I/O solution for extreme-scale HPC will need to manage metadata in enormous volumes, at enormous rates and scales, and, again, in a compromised environment. Current storage and I/O solutions, and even near term designs, do not have the performance or scalability needed to cover this critical area. Further, non-POSIX access methods that could be proposed may imply the need for new types of metadata. Metadata targets are described in the environments table above.

A3-2.3 Data

Future extreme-scale storage systems will need to manage data in enormous volumes, at enormous rates and scales, and in an environment that assumes failure. Billion-way concurrency is expected. Current storage and I/O solutions, and even near term designs, do not have all the attributes needed address future challenges in performance, scalability, error handling, and

concurrency management to name a few. Data performance and concurrency targets are described in the environment table above.

A3-2.4 Quality of Service

Much research has been done on methods for enabling Quality of Service (QOS) for HPC environments, but little of that research has been implemented in products. Future systems might have as few as one job and as many as dozens of jobs running concurrently. Dealing with varying priorities, size, and shape of workloads is very important and such solutions could be very disruptive to existing storage and I/O solutions. Understanding when the complete I/O solution is achieving full potential or struggling in some aspect, for instance bandwidth or small I/Os, would be extremely useful in dealing with simultaneous workloads. Innovative solutions to this QOS problem will be needed, especially fundamental solutions that are not an afterthought. The environment table above describes the concurrency levels and attributes of availability/reliability that are relevant to addressing QOS issues.

A3-2.5 Security

Adequate mechanisms to enforce “need-to-know,” the ability to separate access of data between users and groups of users, are needed. The tension between providing convenient file sharing and proper security remains a challenging problem. Knowing user activity (for example, audit logs) is a desired feature that rarely exists in today’s storage systems. Future storage and I/O solutions for extreme-scale HPC may preclude current solutions entirely. Models that assume complete trust of the clients, relying only on local methods of authentication and authorization at the client, may be inadequate going forward. The data and metadata transaction rates that a scalable security solution must support are described in the environments table above.

A3-2.6 New Device/Topology Exploration/Exploitation

The inability to solve the entire storage and I/O problems in extreme-scale HPC with traditional disk storage devices alone gives rise to the need to develop or exploit new storage and/or network technologies. Current software stacks could exploit new storage devices, networks, or topologies, while completely new I/O stacks together with new storage environments could fundamentally change this entire storage and I/O area for extreme-scale HPC.

For further background information on needed research for extreme-scale HPC storage and I/O, please refer to <http://institute.lanl.gov/hec-fsio/docs/>.

A3-3 Performance Metrics (MR)

Offeror shall estimate or quantify the impact of the proposed technology over industry roadmaps and trends. This information shall be provided for all of the metrics listed below. If Offeror determines that a particular metric is not applicable to the technology being proposed, then Offeror shall explain why they believe the metric is not relevant and shall replace that metric with an alternate *meaningful* metric.

Quantities specified should reflect solutions that are productized in the 2020 timeframe. These metrics are independent, but a solution that can deliver advances in more than one metric is more

desirable than one that solves only one metric at the expense of the others. The most meritorious improvements will make substantial gains over industry roadmaps/trends and substantiate a convincing path to achieving the extreme-scale technology characteristics required by DOE.

Mean Time to Application Visible Interrupt. Offeror shall specify the projected mean time to application visible interrupt caused by the storage system, measured at full system job size.

Peak Burst I/O Rate. Offeror shall specify the projected peak burst I/O rate of the proposed solution. This is defined as the maximum transfer rate to the closest buffer cache.

Data Rate for Unaligned/Variable-Sized Requests. Offeror should specify the projected data rate at which the storage system is able to read and write two system memories having irregular and unaligned data patterns in parallel from 1 billion processes.

Metadata Transaction Rates. Offeror shall specify the projected metadata transaction rates of the proposed solution. Rates for both metadata insertions and queries shall be provided. In this context, storage system metadata includes both the well-known attributes tracked by the storage system itself and the addressable content of any user-extensible spaces it might maintain for the purpose of augmenting a related collection of data.

Metadata Performance Efficiency. Offeror should specify the projected amount of degradation, if any, on metadata query/retrieval operations during storage system peak read and/or write operations.

End-to-End Data Protection. Offeror should describe how end-to-end data protection will be accomplished. Discuss any projected impact on performance and/or functionality of the storage system with end-to-end data protection fully engaged.

A3-4 Target Requirements

The requirements below apply to supercomputers that will be deployed at the end of this decade to meet DOE mission needs. As previously stated, Offerors need not address all problem areas, and thus the Offeror need not respond to TR below if the proposed capability does not address that problem area. In all TR responses that are provided, Offeror should discuss what progress will be made in the next two years and describe what follow-on efforts will be needed to fully achieve these goals. Offeror should describe in detail how the metric will be evaluated, including the measurement method that will be used for example, simulation or prototype) and any assumptions that will be made.

The targets in this section assume the “baseline” system attributes described in the 2020 column of the table above.

A3-4.1 Reliability/Availability/Manageability

Mean Time to Application Visible Interrupt (TR-1)

The mean time to application visible interrupt caused by the storage system measured at full system job size should be no less than 30 days.

Mean Time to Data Loss (TR-1)

The mean time to unrecoverable data loss caused by the storage system should be no less than 120 days, and all data lost can be enumerated by name for system users. Calculation of this metric should assume a storage system that is 80-percent full and continuously performing full supercomputer memory dumps each hour.

Partial Unavailability (TR-1)

Mean time to partial unavailability for the storage system should be no less than 20 days.

Total Unavailability (TR-1)

Mean time to total unavailability for the storage system should be no less than 120 days.

End-to-End Data Integrity with Low Overhead (TR-1)

End-to-end data integrity capability from application interface and back should be provided with no more than 10-percent overhead on metadata insert/query and data movement rates, measured on a full supercomputer-system-sized workload.

A3-4.2 Metadata

Improved Metadata Insert Rates (TR-1)

Transactionally secure insert rates into metadata store with consistency provided in less than 10 s should be no less than 1 million/s.

Significantly Improved Metadata Insert Rates (TR-2)

Transactionally secure insert rates into metadata store with consistency provided in less than 10 s should be no less than 10 million/s.

Greatly Improved Metadata Insert Rates (TR-3)

Transactionally secure insert rates into metadata store with consistency provided in less than 10 s should be no less than 100 million/s.

Improved Metadata Query Rates (TR-1)

Keyed lookup and retrieval of metadata entries should be no less than 100 thousand/s.

Significantly Improved Metadata Query Rates (TR-2)

Keyed lookup and retrieval of metadata entries should be no less than 1 million/s.

Greatly Improved Metadata Query Rates (TR-3)

Keyed lookup and retrieval of metadata entries should be no less than 10 million/s.

Metadata Richness (TR-2)

The storage system should provide the capability for users to annotate data and find data via multiple metadata approaches (for example, hierarchies or key values).

A3-4.3 Data

Improved Data Rates for Unaligned/Variable-Sized Requests (TR-1)

The storage system should be able to read and write two system memories having irregular and unaligned data patterns in parallel from 1 billion processes at 100 TiB/s.

Greatly Improved Data Rates for Unaligned/Variable-sized Requests (TR-2)

The storage system should be able to read and write two system memories having irregular and unaligned data patterns in parallel from 1 billion processes at 300 TiB/s

Improved Data Rates for Many Unaligned/Variable-sized Requests (TR-1)

The storage system should be able to read and write 20 system memories having irregular and unaligned data patterns in parallel at 20 TiB/s.

A3-4.4 QoS

Efficient Metadata Requests During Large Data Movement (TR-1)

The storage system should have no more than 25-percent degradation on metadata query/retrieval operations during storage system peak read and/or write operations.

Highly Efficient Metadata Requests During Large Data Movement (TR-2)

The storage system should have no more than 10-percent degradation on metadata query/retrieval operations during storage system peak read and/or write operations.

Efficient Multiple Concurrent Large Data Movement (TR-1)

The storage system should allow each of four parallel concurrent read/write workloads occupying the entire supercomputer to operate at 75 percent of the data rate these workloads would receive without the other concurrent workloads.

Highly Efficient Multiple Concurrent Large Data Movement (TR-2)

The storage system should allow each of four parallel concurrent read/write workloads occupying the entire supercomputer to operate at 90 percent of the data rate these workloads would receive without the other concurrent workloads.

A3-4.5 Security

End-to-End Data Protection (TR-1)

End-to-end data security capability should be provided, from application interface to storage system and back.

Minimal End-to-End Data Protection Overhead (TR-2)

End-to-end data security capability should be provided, from application interface to storage system and back, while meeting all other TR-1 and TR-2 in this section.